# A grammar compressor for collections of reads with applications to the construction of the BWT

**Diego Díaz-Domínguez**    Gonzalo Navarro

Department of Computer Science
University of Chile

Centre for Biotechnology and Bioengineering (CeBiB)
University of Chile

Data Compression Conference, March 2021

# Motivation

In Genomics, *sequencing reads* are massive and repetitive string collections of raw DNA sequences
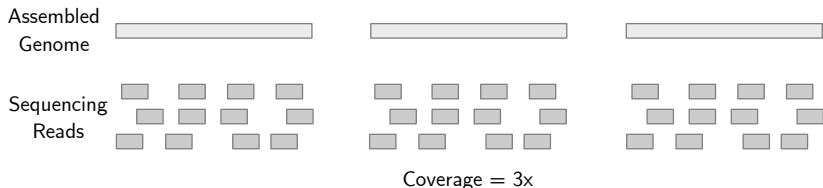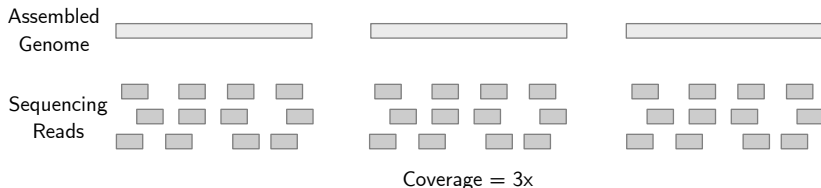
# Motivation

In Genomics, *sequencing reads* are massive and repetitive string collections of raw DNA sequences

# Motivation

In Genomics, *sequencing reads* are massive and repetitive string collections of raw DNA sequences



Coverage $= 3x$

# Motivation

In Genomics, *sequencing reads* are massive and repetitive string collections of raw DNA sequences



Read collections can be much more massive than assembled genomes

# Motivation

## Long-term goal

Producing a compact representation for reads for storage and analysis purposes

## Long-term goal

Producing a compact representation for reads for storage and analysis purposes

Which compression method should we use?

# Motivation

## Long-term goal

Producing a compact representation for reads for storage and analysis purposes

Which compression method should we use?

- Lempel-Ziv (LZ) and Grammars achieve good compression ratios

# Motivation

## Long-term goal

Producing a compact representation for reads for storage and analysis purposes

Which compression method should we use?

- Lempel-Ziv (LZ) and Grammars achieve good compression ratios
  - Performing any analysis on the reads requires decompression

# Motivation

## Long-term goal

Producing a compact representation for reads for storage and analysis purposes

Which compression method should we use?

- Lempel-Ziv (LZ) and Grammars achieve good compression ratios
  - Performing any analysis on the reads requires decompression
  - Lempel-Ziv compresses better, but grammars allow random decompression

# Motivation

## Long-term goal

Producing a compact representation for reads for storage and analysis purposes

Which compression method should we use?

- Lempel-Ziv (LZ) and Grammars achieve good compression ratios
  - Performing any analysis on the reads requires decompression
  - Lempel-Ziv compresses better, but grammars allow random decompression
- The compression ratio of the BWT is not as good as that of LZ or Grammars. However, it is good for indexing

# Motivation

## Long-term goal

Producing a compact representation for reads for storage and analysis purposes

Which compression method should we use?

- Lempel-Ziv (LZ) and Grammars achieve good compression ratios
  - Performing any analysis on the reads requires decompression
  - Lempel-Ziv compresses better, but grammars allow random decompression
- The compression ratio of the BWT is not as good as that of LZ or Grammars. However, it is good for indexing

## What do we propose in this work?

A grammar compressor from which we can efficiently compute the eBWT of the reads

# Related concepts

Induced suffix sorting (ISS) (**Nong _et al_. 2009**) is a technique for building the **SA** of a text $T[1..n]$ in linear time

# Related concepts

Induced suffix sorting (ISS) (**Nong *et al.* 2009**) is a technique for building the **SA** of a text $T[1..n]$ in linear time

## Definitions: (**Nong *et al.* 2009**)

- **L-type** (L): $T[i..n] >_{lex} T[i+1..n]$
- **S-type** (S): $T[i..n] <_{lex} T[i+1..n]$
- **LMS-type** (S*): $T[i-1..n] >_{lex} T[i..n] <_{lex} T[i+1..n]$

A substirng $P = T[i..j]$ is a **LMS-substring** if the suffixes $T[i..n]$ and $T[j..n]$ are S* and no other suffix in $P$ is S*

- Let $\mathcal{T} = \{T_1, T_2, .., T_t\}$ be a collection of $t$ strings with average length $k$

# Definitions

- Let $\mathcal{T} = \{T_1, T_2, .., T_t\}$ be a collection of $t$ strings with average length $k$
- Let $T = T_1\$T_2...\$T_t$ be the concatenation of the elements in $\mathcal{T}$, separated by a dummy symbol \$

# Definitions

- Let $\mathcal{T} = \{T_1, T_2, .., T_t\}$ be a collection of $t$ strings with average length $k$
- Let $T = T_1\$T_2...\$T_t$ be the concatenation of the elements in $\mathcal{T}$, separated by a dummy symbol $
  - The total size of $T$ is denoted as $n$

# Definitions

- Let $\mathcal{T} = \{T_1, T_2, .., T_t\}$ be a collection of $t$ strings with average length $k$
- Let $T = T_1\$T_2...\$T_t$ be the concatenation of the elements in $\mathcal{T}$, separated by a dummy symbol \$
  - The total size of $T$ is denoted as $n$
- Let $\mathcal{G} = \{\Sigma, V, \mathcal{R}, S\}$ be a context-free grammar that only produces $T$

# Definitions

- Let $\mathcal{T} = \{T_1, T_2, .., T_t\}$ be a collection of $t$ strings with average length $k$
- Let $T = T_1\$T_2...\$T_t$ be the concatenation of the elements in $\mathcal{T}$, separated by a dummy symbol \$
  - The total size of $T$ is denoted as $n$
- Let $\mathcal{G} = \{\Sigma, V, \mathcal{R}, S\}$ be a context-free grammar that only produces $T$
  - $\mathcal{R}$ is the set of rules

# Definitions

- Let $\mathcal{T} = \{T_1, T_2, .., T_t\}$ be a collection of $t$ strings with average length $k$
- Let $T = T_1\$T_2...\$T_t$ be the concatenation of the elements in $\mathcal{T}$, separated by a dummy symbol $\$$
    - The total size of $T$ is denoted as $n$
- Let $\mathcal{G} = \{\Sigma, V, \mathcal{R}, S\}$ be a context-free grammar that only produces $T$
    - $\mathcal{R}$ is the set of rules
    - $S$ is the start symbol

# Definitions

- Let $\mathcal{T} = \{T_1, T_2, .., T_t\}$ be a collection of $t$ strings with average length $k$
- Let $T = T_1\$T_2...\$T_t$ be the concatenation of the elements in $\mathcal{T}$, separated by a dummy symbol $\$$
    - The total size of $T$ is denoted as $n$
- Let $\mathcal{G} = \{\Sigma, V, \mathcal{R}, S\}$ be a context-free grammar that only produces $T$
    - $\mathcal{R}$ is the set of rules
    - $S$ is the start symbol
    - $r = |\mathcal{R}|$ is the number of rules

# Definitions

- Let $\mathcal{T} = \{T_1, T_2, .., T_t\}$ be a collection of $t$ strings with average length $k$
- Let $T = T_1\$T_2...\$T_t$ be the concatenation of the elements in $\mathcal{T}$, separated by a dummy symbol \$
    - The total size of $T$ is denoted as $n$
- Let $\mathcal{G} = \{\Sigma, V, \mathcal{R}, S\}$ be a context-free grammar that only produces $T$
    - $\mathcal{R}$ is the set of rules
    - $S$ is the start symbol
    - $r = |\mathcal{R}|$ is the number of rules
    - $g$ is the sum of the lengths of the right hands of the rules

# LMSg algorithm

$T^1 = $ g t a t t a c c $ c t a a t a g t a c c $ g a c c a g a c c a g t $

# LMSg algorithm

$T^1 =$ g t a t t a c c \$ c t a a t a g t a c c \$ g a c c a g a c c a g t \$

    S  L  S*  L  L  S*  L  L  S*  S  L  S*  S  L  S*  S  L  S*  L  L  S*  L  S*  S  L  S*  L  S*  L  L  S*  S  L  S*

# LMSg algorithm

$T^1 =$ g t a t t a c c $ c t a a t a g t a c c $ g a c c a g a c c a g t $

S L S* L L S* L L S* S L S* S L S* S L S* L L S* L S* S L S* L S* L L S* S L S*

# LMSg algorithm

$T^2 =$ 7    **8**    2    **4**    **1**    7    2    5    3    5    3    **6**
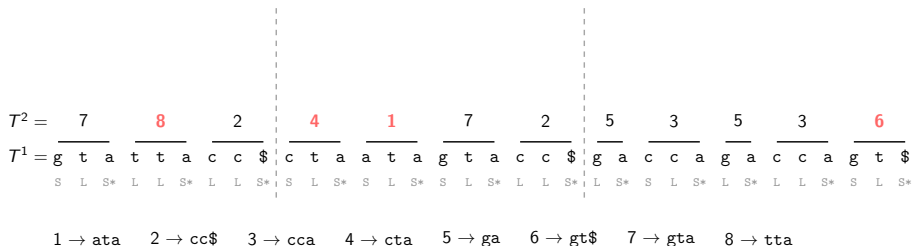
$T^1 =$ g t a t t a c c $ c t a a t a g t a c c $ g a c c a g a c c a g t $

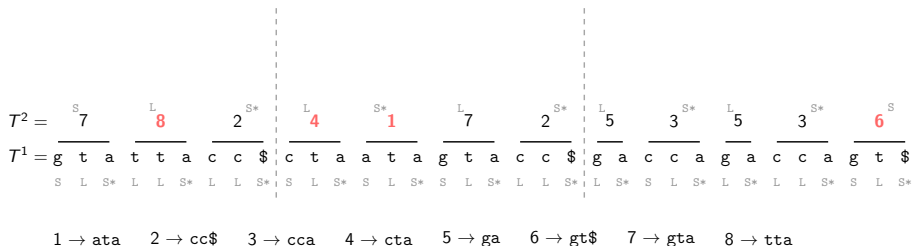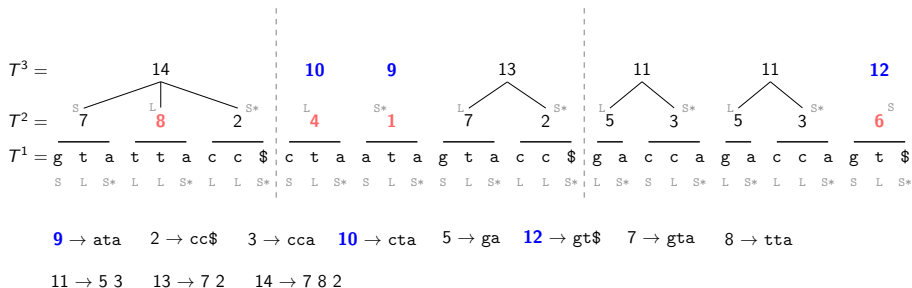S L S* L L S* L L S* S L S* S L S* S L S* L L S* L S* S L S* L S* L L S* S L S*

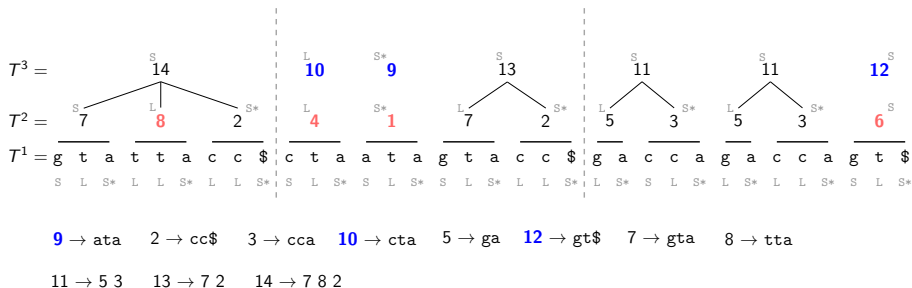$1 \rightarrow$ ata    $2 \rightarrow$ cc$    $3 \rightarrow$ cca    $4 \rightarrow$ cta    $5 \rightarrow$ ga    $6 \rightarrow$ gt$    $7 \rightarrow$ gta    $8 \rightarrow$ tta

# LMSg algorithm

$$T^2 = \overset{\text{S}}{7} \quad \overset{\text{L}}{\mathbf{8}} \quad \overset{\text{S*}}{2} \quad \overset{\text{L}}{\mathbf{4}} \quad \overset{\text{S*}}{\mathbf{1}} \quad \overset{\text{L}}{7} \quad \overset{\text{S*}}{2} \quad \overset{\text{L}}{5} \quad \overset{\text{S*}}{3} \quad \overset{\text{L}}{5} \quad \overset{\text{S*}}{3} \quad \overset{\text{S}}{\mathbf{6}}$$

$$T^1 = \underset{\text{S L S*}}{\text{g t a}} \ \underset{\text{L L S*}}{\text{t t a}} \ \underset{\text{L L S*}}{\text{c c \$}} \ \underset{\text{S L S*}}{\text{c t a}} \ \underset{\text{S L S*}}{\text{a t a}} \ \underset{\text{S L S*}}{\text{g t a}} \ \underset{\text{L L S*}}{\text{c c \$}} \ \underset{\text{L S*}}{\text{g a}} \ \underset{\text{S L S*}}{\text{c c a}} \ \underset{\text{L S*}}{\text{g a}} \ \underset{\text{L L S*}}{\text{c c a}} \ \underset{\text{S L S*}}{\text{g t \$}}$$

$1 \rightarrow \texttt{ata}$  $2 \rightarrow \texttt{cc\$}$  $3 \rightarrow \texttt{cca}$  $4 \rightarrow \texttt{cta}$  $5 \rightarrow \texttt{ga}$  $6 \rightarrow \texttt{gt\$}$  $7 \rightarrow \texttt{gta}$  $8 \rightarrow \texttt{tta}$

# LMSg algorithm

$T^3 =$   14   **10**   **9**   13   11   11   **12**

$T^2 =$   7   **8**   2   **4**   **1**   7   2   5   3   5   3   **6**

$T^1 =$ g t a t t a c c \$ c t a a t a g t a c c \$ g a c c a g a c c a g t \$

**9** → ata    2 → cc\$    3 → cca    **10** → cta    5 → ga    **12** → gt\$    7 → gta    8 → tta

11 → 5 3    13 → 7 2    14 → 7 8 2

# LMSg algorithm

$T^3 =$

$T^2 =$

$T^1 =$ g t a t t a c c \$ c t a a t a g t a c c \$ g a c c a g a c c a g t \$

14

7 8 2 10 9 13 11 11 12

4 1 7 2 5 3 5 3 6

S L S* L L S* L L S* S L S* S L S* S L S* S L S* L L S* L S* S L S* L S* L L S* S L S*

$9 \rightarrow$ ata    $2 \rightarrow$ cc\$    $3 \rightarrow$ cca    $10 \rightarrow$ cta    $5 \rightarrow$ ga    $12 \rightarrow$ gt\$    $7 \rightarrow$ gta    $8 \rightarrow$ tta

$11 \rightarrow 5\ 3$    $13 \rightarrow 7\ 2$    $14 \rightarrow 7\ 8\ 2$

# LMSg algorithm



$T^4 =$     **19**     **16**     **15**     **18**     17

$T^3 =$     14     **10**     **9**     13     11     11     **12**

$T^2 =$     7   **8**   2    **4**    **1**    7   2    5   3    5   3    **6**

$T^1 =$ g t a t t a c c \$ c t a a t a g t a c c \$ g a c c a g a c c a g t \$

$15 \to$ ata    $2 \to$ cc\$    $3 \to$ cca    $16 \to$ cta    $5 \to$ ga    $12 \to$ gt\$    $7 \to$ gta    $8 \to$ tta

$11 \to 5\ 3$    $13 \to 7\ 2$    $14 \to 7\ 8\ 2$    $17 \to 11\ 11\ 12$

# LMSg algorithm



$T^4 =$      **19**      **16**   **15**     **18**        17

$T^3 =$    14     **10**   **9**     13     11    11    **12**

$T^2 =$   7    **8**    2     **4**     **1**     7    2    5    3    5    3    **6**

$T^1 =$ g t a t t a c c \$ c t a a t a g t a c c \$ g a c c a g a c c a g t \$

$\mathcal{R} \begin{cases} \textbf{15} \to \texttt{ata} \quad 2 \to \texttt{cc\$} \quad 3 \to \texttt{cca} \quad \textbf{16} \to \texttt{cta} \quad 5 \to \texttt{ga} \quad \textbf{12} \to \texttt{gt\$} \quad 7 \to \texttt{gta} \quad 8 \to \texttt{tta} \\ 11 \to 5\ 3 \quad 13 \to 7\ 2 \quad 14 \to 7\ 8\ 2 \quad 17 \to 11\ 11\ 12 \end{cases}$

$$S \to 19\ 16\ 15\ 18\ 17$$

# LMSg algorithm



$T^4 =$         **19**      **16**    **15**     **18**            17

$T^3 =$     14    **10**   **9**    13     11    11    **12**

$T^2 =$    7    **8**    2    **4**    **1**    7    2    5    3    5    3    **6**

$T^1 =$ g t a t t a c c \$ c t a a t a g t a c c \$ g a c c a g a c c a g t \$

$$\mathcal{R} \begin{cases} \textbf{15} \rightarrow \texttt{ata} \quad 2 \rightarrow \texttt{cc\$} \quad 3 \rightarrow \texttt{cca} \quad \textbf{16} \rightarrow \texttt{cta} \quad 5 \rightarrow \texttt{ga} \quad \textbf{12} \rightarrow \texttt{gt\$} \quad 7 \rightarrow \texttt{gta} \quad 8 \rightarrow \texttt{tta} \\ 11 \rightarrow \texttt{5 3} \quad 13 \rightarrow \texttt{7 2} \quad 14 \rightarrow \texttt{7 8 2} \quad 17 \rightarrow \texttt{11 11 12} \end{cases}$$

$$S \rightarrow 19\ 16\ 15\ 18\ 17$$

## Time complexity

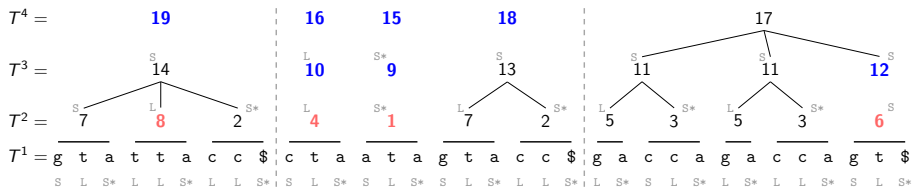The grammar construction takes $O(n \log k)$ time

# LMSg algorithm



$T^4 = $  **19**   **16**   **15**   **18**   17

$T^3 = $  14   **10**   **9**   13   11   11   **12**

$T^2 = $  7   **8**   2   **4**   **1**   7   2   5   3   5   3   **6**

$T^1 = $ g t a t t a c c $ c t a a t a g t a c c $ g a c c a g a c c a g t $

$BWT(T^3)$

| | |
|---|---|
| 10 | 9 |
| 13 | 10 |
| 12 | 11 |
| 11 | 11 |
| 11 | 12 |
| 9 | 13 |
| 14 | 14 |

# LMSg algorithm



We call the replacement of a nonterminal $BWT(T^i)[j]$ its *partial decompression*

1 4 3 2

5 3 2 0

     8 1

2 2 3 2

3 3 12 8

  1 3 2

A range of partially decompressed phrases in some $BWT(T^i)$

```
1 4  3 2
5 3  2 0
     8 1
2 2  3 2
3 3 12 8
   1 3 2
```

A range of partially decompressed phrases in some $BWT(T^i)$

Strings 3 2 and 3 2 0 are two distinct suffixes in partially decompressed phrases

# Sketch for inferring the eBWT

| 1 4 | 3 2 | 5 | 3 2 0 L L S |
|-----|-----|---|-------------|
| 5 3 | 2 0 | 4 | 3 2   L S |
|     | 8 1 | 2 | 3 2 |
| 2 2 | 3 2 | 1 | 3 2 |
| 3 3 | 12 8 | | |
| 1 | 3 2 | | |

# Sketch for inferring the eBWT

```
1 4  3 2    5   3 2 0 L L S

5 3 2 0    4   3 2   L S

     8 1    2   3 2

2 2  3 2    1   3 2

3 3 12 8

   1  3 2
```

In one scan of $BWT(T^i)$, we obtain the left contexts of a specific phrase's suffix

# Sketch for inferring the eBWT

| 1 | 4 | 3 | 2 | | 5 | 3 | 2 | 0 | L | L | S |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 5 | 3 | 2 | 0 | | 4 | 3 | 2 | | L | S | |
| | | 8 | 1 | | 2 | 3 | 2 | | | | |
| 2 | 2 | 3 | 2 | | 1 | 3 | 2 | | | | |
| 3 | 3 | 12 | 8 | | | | | | | | |
| | 1 | 3 | 2 | | | | | | | | |

In one scan of $BWT(T^i)$, we obtain the left contexts of a specific phrase's suffix

During the construction of $\mathcal{G}$, we encapsulate the repeated suffixes in new rules

```
1 4  3 2      5    3 2 0 L L S

5 3 2 0       4    3 2    L S

     8 1      2    3 2

2 2  3 2      1    3 2

3 3 12 8

   1  3 2
```

In one scan of $BWT(T^i)$, we obtain the left contexts of a specific phrase's suffix

During the construction of $\mathcal{G}$, we encapsulate the repeated suffixes in new rules

## Observation

Building $BWT(T^{i-1})$ reduces *mostly* to sort the *distinct* suffixes in the partial decompressions

# Sketch for inferring the eBWT

|   |   |    |   |   |   |   |   |
|---|---|----|---|---|---|---|---|
| 1 | 4 | 3  | 2 | S | L | L | S |
| 5 | 3 | 2  | 0 | L | L | L | S |
|   |   | 8  | 1 |   |   | L | S |
| 2 | 2 | 3  | 2 | S | S | L | S |
| 3 | 3 | 12 | 8 | S | S | L | S |
|   | 1 | 3  | 2 |   | S | L | S |

We do not have enough information for sorting suffixes of length 1

We use $BWT(T^i)$ to obtain the right context of suffix 1 in row 3

## Observation

Building $BWT(T^{i-1})$ reduces *mostly* to sort the *distinct* suffixes in the partial decompressions

# Representing the grammar

We adapt the **grammar tree** data structure proposed by **Claude *et al.* 2012**

We adapt the **grammar tree** data structure proposed by **Claude *et al.* 2012**
**Procedure**:

# Representing the grammar

We adapt the **grammar tree** data structure proposed by **Claude *et al.* 2012**
**Procedure**:

- We traverse the parse tree of $\mathcal{G}$ in level-order

# Representing the grammar

We adapt the **grammar tree** data structure proposed by **Claude *et al.* 2012**
**Procedure**:

- We traverse the parse tree of $\mathcal{G}$ in level-order
- Terminals are stored as leaves

# Representing the grammar

We adapt the **grammar tree** data structure proposed by **Claude et al. 2012**
**Procedure**:

- We traverse the parse tree of $\mathcal{G}$ in level-order
- Terminals are stored as leaves
- The first occurrence of a nonterminal is stored as an internal node

# Representing the grammar

We adapt the **grammar tree** data structure proposed by **Claude *et al.* 2012**
**Procedure**:

- We traverse the parse tree of $\mathcal{G}$ in level-order
- Terminals are stored as leaves
- The first occurrence of a nonterminal is stored as an internal node
    - The next occurrences are stored as leaves

# Representing the grammar

We adapt the **grammar tree** data structure proposed by **Claude *et al.* 2012**
**Procedure**:

- We traverse the parse tree of $\mathcal{G}$ in level-order
- Terminals are stored as leaves
- The first occurrence of a nonterminal is stored as an internal node
  - The next occurrences are stored as leaves
  - The leaf's label is the first occurrence of the nonterminal

# Representing the grammar

We adapt the **grammar tree** data structure proposed by **Claude *et al.* 2012**
**Procedure**:

- We traverse the parse tree of $\mathcal{G}$ in level-order
- Terminals are stored as leaves
- The first occurrence of a nonterminal is stored as an internal node
  - The next occurrences are stored as leaves
  - The leaf's label is the first occurrence of the nonterminal
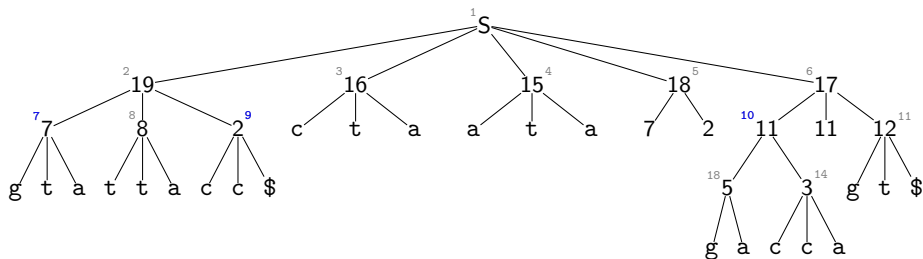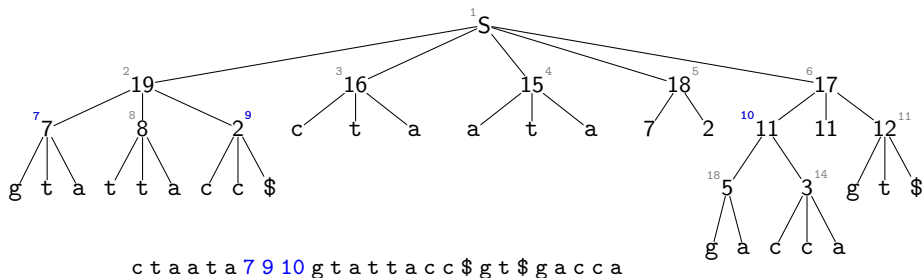- We encode the topology in **LOUDS**

# Representing the grammar

We adapt the **grammar tree** data structure proposed by **Claude *et al.* 2012**
**Procedure**:

- We traverse the parse tree of $\mathcal{G}$ in level-order
- Terminals are stored as leaves
- The first occurrence of a nonterminal is stored as an internal node
    - The next occurrences are stored as leaves
    - The leaf's label is the first occurrence of the nonterminal
- We encode the topology in **LOUDS**
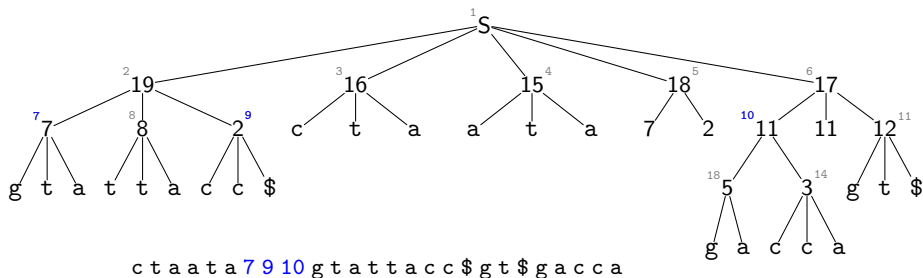- The leaf labels are stored in a succinct array

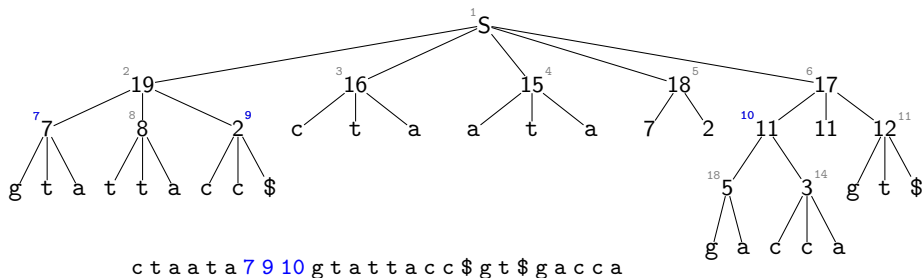# Representing the grammar

# Representing the grammar

# Representing the grammar



## Space usage

The grammar tree uses $2g + (g - r)\log(r)$ bits of space

# Representing the grammar



## Space usage

The grammar tree uses $2g + (g - r)\log(r)$ bits of space

We have to reconstruct the original nonterminal symbols to get eBWT

# Experiments

We compressed read collections from five human samples of the Human Genome Diversity Project

# Experiments

We compressed read collections from five human samples of the Human Genome Diversity Project

| Number of genomes | Uncompressed size (GB) |
|:---:|:---:|
| 1 | 12.77 |
| 2 | 23.43 |
| 3 | 34.30 |
| 4 | 45.89 |
| 5 | 57.37 |

# Experiments

We compressed read collections from five human samples of the Human Genome Diversity Project

| Number of genomes | Uncompressed size (GB) |
| --- | --- |
| 1 | 12.77 |
| 2 | 23.43 |
| 3 | 34.30 |
| 4 | 45.89 |
| 5 | 57.37 |

We evaluated the following methods:

# Experiments

We compressed read collections from five human samples of the Human Genome Diversity Project

| Number of genomes | Uncompressed size (GB) |
|:---:|:---:|
| 1 | 12.77 |
| 2 | 23.43 |
| 3 | 34.30 |
| 4 | 45.89 |
| 5 | 57.37 |

We evaluated the following methods:

- LMSg (LPG)

# Experiments

We compressed read collections from five human samples of the Human Genome Diversity Project

| Number of genomes | Uncompressed size (GB) |
|:---:|:---:|
| 1 | 12.77 |
| 2 | 23.43 |
| 3 | 34.30 |
| 4 | 45.89 |
| 5 | 57.37 |

We evaluated the following methods:

- LMSg (LPG)
- RePair + Prefix-Free Parsing (BR)

# Experiments

We compressed read collections from five human samples of the Human Genome Diversity Project

| Number of genomes | Uncompressed size (GB) |
|:---:|:---:|
| 1 | 12.77 |
| 2 | 23.43 |
| 3 | 34.30 |
| 4 | 45.89 |
| 5 | 57.37 |

We evaluated the following methods:

- LMSg (LPG)
- RePair + Prefix-Free Parsing (BR)
- p7zip (highest compression mode) (7Z)

# Experiments

We compressed read collections from five human samples of the Human Genome Diversity Project

| Number of genomes | Uncompressed size (GB) |
|:---:|:---:|
| 1 | 12.77 |
| 2 | 23.43 |
| 3 | 34.30 |
| 4 | 45.89 |
| 5 | 57.37 |

We evaluated the following methods:

- LMSg (LPG)
- RePair + Prefix-Free Parsing (BR)
- p7zip (highest compression mode) (7Z)
- FM-index (without *SA*) (FM)

# Experiments

We compressed read collections from five human samples of the Human Genome Diversity Project

| Number of genomes | Uncompressed size (GB) |
| :---: | :---: |
| 1 | 12.77 |
| 2 | 23.43 |
| 3 | 34.30 |
| 4 | 45.89 |
| 5 | 57.37 |

We evaluated the following methods:

- LMSg (LPG)
- RePair + Prefix-Free Parsing (BR)
- p7zip (highest compression mode) (7Z)
- FM-index (without *SA*) (FM)
- RLFM-index (without *SA*) (RLFM)

# Experiments

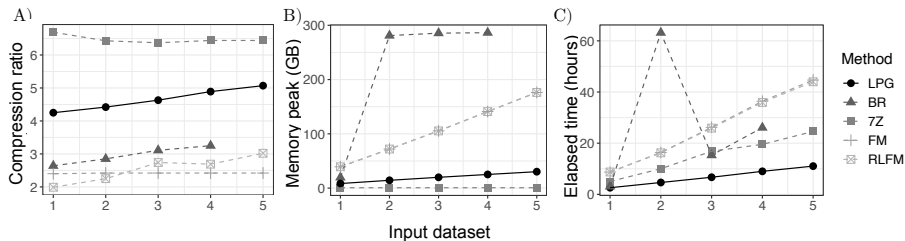We compressed read collections from five human samples of the Human Genome Diversity Project

| Number of genomes | Uncompressed size (GB) |
|:---:|:---:|
| 1 | 12.77 |
| 2 | 23.43 |
| 3 | 34.30 |
| 4 | 45.89 |
| 5 | 57.37 |

We evaluated the following methods:

- LMSg (LPG)
- RePair + Prefix-Free Parsing (BR)
- p7zip (highest compression mode) (7Z)
- FM-index (without *SA*) (FM)
- RLFM-index (without *SA*) (RLFM)

    We use 10 threads with all the methods (when possible)

# Results



A) Compression ratio — Input dataset (1–5)
B) Memory peak (GB) — Input dataset (1–5)
C) Elapsed time (hours) — Input dataset (1–5)

Method
- LPG
- BR
- 7Z
- FM
- RLFM

The compression ratio was measured as the size of the plain representation divided by the compressed representation

- Simplify the grammar and apply RePair on top of it

# Further work

- Simplify the grammar and apply RePair on top of it
- Remove dollar symbols

# Further work

- Simplify the grammar and apply RePair on top of it
- Remove dollar symbols
- A better implementation of the parallel decompression

# Further work

- Simplify the grammar and apply RePair on top of it
- Remove dollar symbols
- A better implementation of the parallel decompression
- Create a Self-Index from the grammar

# Further work

- Simplify the grammar and apply RePair on top of it
- Remove dollar symbols
- A better implementation of the parallel decompression
- Create a Self-Index from the grammar
- Build other data structures: LCP, de Bruijn graphs ...

# Questions?